# **ReWIRED:** Instructional Explanations in Teacher-Student Dialogues

**Anonymous ACL submission** 

#### Abstract

How to assess the quality of teaching in instruc-001 tional explanation dialogues is a recurring point of debate in didactics research. For the NLP community, this is a challenging topic thus far, even with the use of LLMs. To address the matter, we create a new annotation scheme of 007 teaching acts aligned with contemporary didactic teaching models. On this basis, we extend an existing dataset of conversational explanations about communicating scientific understanding 011 in teacher-student settings on five levels of the explainee's expertise, with the proposed teaching annotation: explanation and dialogue acts. For better granularity, we reframe the task from 015 a dialogue turn classification to a span labeling task. We then evaluate language models on the 017 labeling of such acts and find that the broad range and structure of the proposed labels is hard to model for LLMs such as GPT-3.5/-4 019 via prompting, but a fine-tuned BERT can perform both act classification and span labeling well. Finally, we operationalize a series of quality metrics for instructional explanations in the form of a test suite. We find that they match the five expertise levels well and that experts in our data often stick to best practices in teaching.

### 1 Introduction

027

037

041

The recent paradigm shift in NLP towards LLMs such as ChatGPT has impacted cross-disciplinary research with education and other social sciences. However, automating teacher coaching (Wang and Demszky, 2023) and student tutoring (Macina et al., 2023) has shown limited success so far. A recent work in tutoring by Lee et al. (2023) has explored creating interactive dialogues to answer children's why and how questions. Measures for estimating the quality of discourse (McNamara et al., 2014) or model-generated explanations (Schuff et al., 2023) exist, but it is unclear how we can assess the quality of teaching in such instructional explanation dialogues and also consider the expertise level of the



Figure 1: Instructional explanation dialogue of an expert (center) explaining machine learning to a child (left). Labels on the right indicate the teaching act associated with the turn(s) or span(s) with the same color.

#### explainee (Wachsmuth and Alshomary, 2022).

In this work, we first propose a scheme of teaching acts that connect dialogical surface-level utterances with the processes described by two popular teaching models (§2). Thereby, we open the doors to the large-scale analysis of teaching strategies, a goal much sought after in didactics (Matsumura et al., 2008). Secondly, we re-annotate the WIRED dataset from Wachsmuth and Alshomary (2022) to include our scheme of teaching acts, expanding on the two act sets from the original, dialogue acts and explanation acts (§3). The dataset is further enhanced by the inclusion of 45 new conversation transcripts, and by a switch from a turn-labeling to a span-labeling setting for higher granularity. We evaluate state-of-the-art language models of different sizes on both turn classification and span labeling (§4) and find that large closed-source models cannot perform either task reasonably well and is easily beaten by a fine-tuned BERT. Lastly, to measure "good teaching" according to didactics research in terms of both *meaning* and *form* (Bender and Koller, 2020), we implement a series of quality metrics for instructional explanations, taking into account the presence and order of teaching acts as well as frequency of explanatory patterns. We dub this new test suite IXQUISITE (§5.3) and find that the metrics correlate well with the five expertise levels in our dataset.

057

061

062

063

071

077

082

090

100

102

103

104

105

106

With the results and findings of this paper, we contribute to both fields, NLP and didactics: To NLP, a more accurate representation of the complex sociolinguistic goals affecting the enactment of effective instructional explanations as well as a sanity check on LLM-based tutoring; to didactics, a new way to look at teaching and lesson-planning at massive scale, by taking a bottom-up approach to modeling the learning and teaching process.

### 2 Background and Related Work

There are many concepts that are common to didactics but are neglected in NLP research. Neither tutoring-related works (Lee et al., 2023; Stasaski et al., 2020) nor concept explanation datasets (Dinan et al., 2019; Jansen et al., 2018) distinguish the type of explanation in social sciences (Miller, 2019) from the interpretation in NLP research.

In science teaching, an explanation is viewed as a practice (or even a purpose) of science or scientists that systematically addresses the questions of "how" and "why" (Kulgemeyer, 2018). Here, instructional explanations are those that aim to "communicate a new cognitive model for understanding the world, or how to perform a task, from one understanding-having interlocutor to an understanding lacking one". While most explainability literature has mostly focused on a more philosophical understanding explanation, as that which connects explanans and explanandum (Miller, 2019), the instructional perspective is closely aligned with the much-needed interest in context for explanations (Mostafazadeh et al., 2020). Despite many systems posing to perform instructional tasks, to our knowledge, they do not take any teaching or learning models into consideration.

Teaching models are frameworks to teach teach-

ers how to plan lessons towards better learning outcomes by structuring lessons in accordance with a psychological model of learning. While there have been attempts at unifying multiple teaching and learning models (explaining how learning happens in the mind of the students) (Oser and Baeriswyl, 2002), many remain skeptical about the feasibility (Allensworth et al., 2008). The actual instantiation of them in real-world classroom environments is affected by many socio-cultural elements (Ball and Rowan, 2004), which make it hard to evaluate teaching at scale (Matsumura et al., 2008) and objectively, without considering other teaching activities and social holistic processes surrounding the explanation (Roelle et al., 2015). Boston (2012) abstracted the differences and used broad definitions of the processes, often leading to positive outcomes, but their approach cannot evaluate low-level, dialogical components of the teaching in a classroom. In this paper, we represent teaching processes (1) in the form of teaching acts (Table 1, Table 5) and investigate if language models can capture the distinctions, and (2) as explanation quality measures (Table 6) and an analysis of how well they correlate with expertise levels of the explainee.

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

Tutoring datasets Our work is closest to Wachsmuth and Alshomary (2022): We re-annotate and extend their dataset, perform similar analyses in terms of statistics and LM experiments, but add a new angle to the data with teaching acts and span-level labeling, allowing us to derive quality events in instructional explanations (§5.3) and experiments with LLMs (§5.2). In contrast to CIMA (Stasaski et al., 2020), TSCC-2 (Caines et al., 2022), and NCTE (Demszky and Hill, 2023), their dataset was a good target for modelling different teaching types, as the varied levels also highlight how teaching can change depending on educational level and course subject. Kupor et al. (2023) annotated instruction talk moves in classroom settings and their LMs could perform this classification task well, whereas Macina et al. (2023) and Wang and Demszky (2023) were less successful for applying similar models in neural dialogue tutoring.

**Evaluation of instructional explanations** Previous work in this direction include COH-METRIX, a related suite of measures to assess the quality and readability in discourse automatically (McNamara et al., 2014). Schuff et al. (2023) have also proposed proxy measures for explanation quality based on syntactic and model-based text genera-

Dialogue acts	Explanation acts	Teaching acts
<b>D01</b> : Check Question	<b>E01</b> : Test Understanding	T01: Assess Prior Knowledge
Asking a check question	Checking whether the listener understood	Checking what the student knows
	what was being explained	before starting a lesson
D02: What/How Question	E02: Test Prior Knowledge	T02: Lesson Proposal
Asking a what or a how question	Checking the listener's prior	Proposing the steps that will be taken during the lesson
	knowledge of the turn's topic	
<b>D03</b> : Other Question	<b>E03</b> : Provide Explanation	T03: Active Experience
Asking any other question	Explaining any concept or topic	Providing the student with puzzle/question to explore;
	to the listener	(Student:) Interacting with a mental concept
<b>D04</b> : Confirming Answer	E04: Request Explanation	T04: Reflection
Answering a question	Requesting any explanation	Finding gaps in knowledge or inconsistencies;
with confirmation	from the listener	Asking questions about the experience or concept
<b>D05</b> : Disconfirming Answer	E05: Signal Understanding	T05: Knowledge Statement
Answering a question	Informing the listener that	Stating the concept(s) being taught via rules or facts
with disconfirmation	their last utterance was understood	
<b>D06</b> : Other Answer	E06: Signal Non-understanding	T06: Comparison
Giving any other answer	Informing the listener that	Considering similarities and differences between
	the utterance was not understood	the main concept and other related topics or facts
<b>D07</b> : Agreeing Statement	E07: Provide Feedback	<b>T07</b> : Generalization
Conveying agreement on the	Responding qualitatively to an	Exploring how the concept applies to new scenarios,
last utterance of the listener	utterance by correcting errors	experiences and situations outside of the lesson topic
<b>D08</b> : Disagreeing Statement	E08: Provide Assessment	T08: Test Understanding
Conveying disagreement on the	Assessing the listener by rephrasing	Finding out if the concept previously established
last utterance of the listener	their utterance or giving a hint	was received correctly and is properly understood
<b>D09</b> : Informing Statement	<b>E09</b> : Provide Extraneous Information	T09: Engagement Management
Providing information with respect	Giving additional information	Maintaining the classroom context to facilitate effective
to the topic stated in the turn	to foster a complete understanding	teaching, creating rapport between teacher and student
<b>D10</b> : Other Act	E10: Other Act	T10: Other Act

Table 1: Dialogue, explanation and teaching acts (alongside descriptions) in our ReWIRED dataset.

#	Торіс	Explainer
14	Memory	Daphna Shohamy
15	Zero-knowledge proofs	Amit Sahai
16	Black holes	Janna Levin
17	Quantum computing	Talia Gershon
18	Quantum sensing	Chandrasekhar
		Ramanathan
19	Fractals	Keenan Crane
20	Internet	Jim Kurose
21	Moravecs Paradox	Chelsea Finn
22	Infinity	Emily Riehl

Table 2: New topics and explainers in ReWIRED on top of the 13 original topics of WIRED which can be found in Wachsmuth and Alshomary (2022).

tion metrics but found low correlation with human judgments. Demszky et al. (2021) develop a framework for measuring teachers' uptake (defined as *building on the student's contribution via, for example, acknowledgement, repetition or elaboration*). Whitehill and LoCasale-Crouch (2024) explore how LLMs can be used to estimate what they define as "instructional support" domain scores with the help of an observation protocol.

## **3** The **ReWIRED** Dataset

158

159

161

162

163

164

165

166

167

168

169

170

171

Wachsmuth and Alshomary (2022) classified parts of instructional explanation dialogues from a dataset collected from the 5-levels video series, in which an expert in a topic, such as black holes, or music harmony, explains the topic to people of<br/>varying expertise levels:1721. Child,1742. Teenager,1753. Undergraduate college student,1764. Graduate student,1775. Colleague (another expert).178

179

180

181

182

183

184

185

187

188

189

190

191

192

193

194

195

196

We can see that a great deal of the content actually comprises instructional explanations, especially the lower we go on the educational scale. Wachsmuth and Alshomary (2022) introduced two types of conversational acts and used them to model explanation dynamics between explainer and explainee.

To increase the models' awareness of teaching perspectives, we add a new scheme of teaching acts to their original two dimensions (Table 1 with supplementary examples in Appendix A) and carry out a refined annotation process. We further improve on their work by switching from a turn labelling task (classification) to a span labeling one for greater granularity, leading to better semantic performance. We dub this improved dataset **ReWIRED**. In the following, we will introduce these teaching acts and our annotation process.

D01	14049	515	658	57	12	19	59	14	674	E01	2331	214	508	459	220	5	155	197	28	T01	9240	142	1365	1104	909	832	268	101	440
D02	539	13458	719	4	0	22	45	0	480	E02	293	5760	388	406	64	12	62	170	233	T02	257	2439	657	188	500	348	307	87	102
D03	719	391	2976	2	17	24	142	0	955	E03	250	384	209109	1056	1325	99	2980	2019	665	T03	513	201	36459	2238		4279	992	150	852
D04	119	50	2	3012	110	232	210	12	619	E04	362	293	854	16875	150	11	226	469	0	T04	772	255	6429	11865		3601	2508	737	1377
D05	31	17	3	99	1119	52	66	0	175	E05	28	117	664	255	7050	42	673	168	224	T05	458	95	1562	1065	49152	3423	841	157	166
D06	0	26	26	207	96	2400	151	0	143	E06	0	9	59	99	125	933	26	39	55	T06	332	24	1749	474	2792	5127	1302	27	174
D07	319	171	98	219	91	182	6771	4	1181	E07-	30	7	1087	333	656	32	5001	619	219	T07	122	109	2773	662	4249	2065		121	607
D08	10	25	0	94	33	9	153	399	254	E08	376	441	729	426	34	0	289	4626	59	T08	110	60	1143	771	679	716	219	3594	128
D09	978	464	316	676	163	1071	881	42	216711	E09-	49	179	979	262	425	6	458	94	9165	T09	154	94	953	530	1360	396	309	93	6825
	002	002 .	003	00 <sup>A</sup>	005	006	001	008	009		<sup>207</sup>	402	403	£0A	405	£06	40 <sup>1</sup>	£08	40 <sup>9</sup>		107	102	10 <sup>3</sup> .	(0 <sup>A</sup>	1 <sup>65</sup>	106	101	108	10 <sup>9</sup>

Figure 2: ReWIRED inter-annotator agreements for the three dimensions dialogue (left), explanation (center) and teaching (right) on token level. For better visibility, we have scale-adjusted the colors by np.log1p(...)<sup>3</sup>. Each cell shows the number of tokens for which annotators (dis)agreed on a label in a pairwise comparison.

is son	mething people have been thinking about for many years. <u>So how has our</u> •T08 - Test U
Teenager: Is there anything that like gets affected on Earth because of those waves?	ersation changed your understanding of what fractals are all about?
•T07 - Genera Under	rgrad: <u>1 think it's really interesting to see the different ways, fractals will be not</u> •T08 - Test U
Explainer: It's a really good question. Only this instrument, and that's why it was so •T07 - Genera only u	useful, but necessary in being able to render these games and these different
hard to build. And by the time it gets here, it's so weak that it's only squeezing and program	rams that are interesting in the metaverse or different media to be really
stretching space at like the fraction of a nucleus over very large distances. Has your	tiful.

Figure 3: Examples for teaching acts T07 (Generalization) and T08 (Test Understanding).

#### 3.1 Teaching acts

197

198

199

201

205

206

207

210

211

212

213

214

216

217

218

219

220

221

Expanding on Wachsmuth and Alshomary (2022), we present a scheme of teaching acts with which to classify dialogues in instructional settings that are coherent with three current and well-accepted teaching models (§2): Teaching as problem solving (**PS**), teaching as concept building (**CB**) (Krabbe et al., 2015), and Oser and Baeriswyl's unified teaching choreographies (**UT**). This is in line with prior work modeling discourse structure in explanations (Bourse and Saint-Dizier, 2012). Concretely, the acts are described in Table 1. Their connection to teaching models and an example<sup>1</sup> are as follows:

• T01:	Assess Prior Knowledge	(CB, UT).
--------	------------------------	-----------

- **T02**: Lesson Proposal (UT).
- **T03**: *Active Experience* (CB, UT).
- **T04**: *Reflection* (PS).
  - **T05**: *Knowledge Statement* (PS).
- **T06**: *Comparison* (UT).
  - **T07**: *Generalization* (CB, PS), e.g. Figure 3.
    - **T08**: *Test Understanding* (CB), e.g. Figure 3.
    - **T09**: Engagement Management .
    - **T10**: *Other Act*: Any other act that does not fit the above nine acts should instead be placed here.

The main goal of the acts is to bring processes

from teaching models closer to the product of their instantiation in actual dialogue (Stolcke et al., 2000), in a way that parts of the dialogue serve as reasonable evidence that the deep processes predicted by teaching models indeed take place.

#### 3.2 Annotation

For our annotation task, we asked nine in-house researchers from a (computational) linguistics background to participate in our annotation study. The total of 110 transcripts from 22 topics across five expertise levels (Table 2) were separated into three groups, such that every annotator group annotated the entire dataset exactly once, one third for dialogue acts, one third for explanation acts (using the original act description by Wachsmuth and Alshomary, 2022), and finally one third for our new set of teaching acts. This yielded three sets of annotations for the entire dataset concerning all three acts. This is first to reduce the possibility of bias, as some acts are very similar and annotators might be tempted to just repeat previous annotations; and second to reduce costs. For our annotation platform, we used DOCCANO (Nakayama et al., 2018), which alleviated the span-labeling task. We additionally randomized all conversations to reduce bias further.

Our inter-annotator agreements are at Fleiss'  $\kappa = 0.83$  (dialogue acts), 0.79 (explanation acts) and 0.46 (teaching acts). We plot the nine main

246

247

248

249

250

222

223

<sup>&</sup>lt;sup>1</sup>Acts with a colored border have an example in both Figure 1 and Figure 8.



Figure 4: Distribution of teaching acts in ReWIRED across the five expertise levels. Dialogue and explanation act distributions are visualized in Appendix B.

labels of each annotation dimension in Figure 2. They show that there also is quite a bit of uncertainty and confusion regarding our teaching acts because our annotators are knowledgeable in computational linguistics but not so much in pedagogy and didactics. Often confused are E03 and E09, as there is a fine line between what we can deem part of an explanation and what is rather supplementary information, and T06 and T07, since both are about "zooming out" of the topic in question and making a broader set of connections to it. The results of our annotation process are visualized via the distribution of teaching acts in Figure 4.

Our annotation scheme differs from DAMSL (Core and Allen, 1997) and ISO 24617-2 (Bunt et al., 2012) in the granularity of annotations. While there are no dependency relations allowing link structures as in the latter, our annotation scheme enables fine-grained annotation of semantics related to teaching models. Most similar to ours is the CMA schema by Del-Bosque-Trevino et al. (2021) for one-to-one tutorial dialogue sessions: In terms of labels, it vaguely mirrors a lot of acts across all three dimensions, but conflates crucial acts (e.g., *FIM* can be either T02 or T09) and misses out on teaching-related concepts such as Active Experience, Reflection, and Generalization.

### 4 Experiments

To evaluate language models on detecting acts across act dimensions, we conduct two experiments: One on turn-level classification, reproducing Wachsmuth and Alshomary (2022), and one on span-labeling for ReWIRED. For both, we test the hypothesis that fine-tuning a masked LM is more consistent at assigning labels on token-level than LLMs prompted for JSON responses indicating spans and labels.

287

289

290

291

292

293

295

296

297

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

328

329

332

**Classifying acts** For the turn-level classification of dialogue and explanation acts provided by the original WIRED data, we choose the following baselines: SVM with linear kernel for multi-class classification based on MiniLM sentence embeddings (Reimers and Gurevych, 2019), and the top-performing BERT from Wachsmuth and Alshomary (2022). We compare the following LMs on both tasks: BERT for turn-level classification, 110M params; Stable Beluga 2 (SB2) (Mahan et al., 2023; Mukherjee et al., 2023), a type of Llama-2 model (Touvron et al., 2023), 70B params; GPT-3.5-turbo-0613. We provide details on the models in Appendix C.

**Sequence-labeling acts** For the token-level span labeling task of the three annotation dimensions (Table 1) in our new **ReWIRED** dataset, we analyze the capabilities of the following LMs: As a baseline, a BERT for token-level classification with 5-fold cross-validation. We compare it to three prompt-based LLMs: Stable Beluga 2; GPT-3.5-turbo-0613; GPT-4-0125-preview. We provide details on the prompt design for the latter three in Appendix D.

### 5 Results and Discussion

#### 5.1 Classifying acts

We show the best performance we were able to attain in automatic act classification for all three acts using several LLMs, and compare our results with the results of Wachsmuth and Alshomary (2022).

Table 3 shows that LLMs perform poorly in turnlevel dialogue act classification, except for capturing disagreeing statements and answers (D08, D05). The fine-tuned BERT model outperforms all other approaches by a substantial amount. This is also repeated for the explanation act classification: LLMs only excel in recognizing signals of (non-)understanding. Across all sets of classes, however, we also find that none of the approaches is able to capture the labels with a very low amount of data points (D05, D08, E01, T02; see Tables 4 & 9).

#### 5.2 Sequence-labeling acts

Our results for span-level act prediction (Table 4) reveal that this task is very challenging for the LLMs, since they were not fine-tuned on the task.

Dialogue acts	D01	D02	D03	D04	D05	D06	D07	D08	D09	D10	Macro-F <sub>1</sub>
W&A BERT-seq	76.00 %	72.00 %	0.00 %	35.00 %	67.00 %	0.00 %	69.00 %	0.00 %	87.00 %	61.00 %	47.00 %
SVM + SentTf	64.30 %	59.55 %	0.00~%	7.14 %	86.96 %	7.69 %	76.28 %	0.00~%	83.30 %	68.57 %	68.71 %
BERT	87.35 %	<b>82.81</b> %	0.00~%	0.00~%	80.77 %	0.00~%	82.04 %	0.00~%	94.62 %	<b>76.77</b> %	81.67 %
SB2	20.00 %	41.51 %	0.00~%	14.29 %	100.00 %	0.00~%	28.57 %	0.00~%	78.67 %	0.00~%	31.45 %
GPT-3.5	14.33 %	43.36 %	4.41 %	19.15 %	37.93 %	5.92 %	21.41 %	8.00 %	69.51 %	33.88 %	25.79 %
Expl. acts	E01	E02	E03	E04	E05	E06	E07	E08	E09	E10	Macro- $F_1$
W&A BERT-seq	27.00 %	64.00 %	84.00 %	64.00 %	33.00 %	21.00 %	60.00 %	15.00 %	8.00 %	56.00 %	43.00 %
SVM + SentTf	6.90 %	66.34 %	81.37 %	37.89 %	13.84 %	0.00~%	72.99 %	0.00~%	<b>28.07</b> %	55.81 %	63.23 %
BERT	0.00 %	<b>73.05</b> %	93.71 %	<b>78.26</b> %	5.52 %	0.00~%	<b>74.89</b> %	0.00~%	0.00~%	66.04 %	66.67 %
SB2	13.79 %	46.60 %	81.63 %	48.89 %	43.53 %	18.18~%	15.13 %	0.00~%	9.68 %	0.00~%	27.74 %
GPT-3.5	16.87 %	38.76 %	71.70~%	23.30 %	37.00 %	28.30 %	5.06 %	0.00~%	2.86 %	27.85 %	27.17 %

Table 3: Language models evaluated on the tasks of classifying dialogue and explanation acts of whole dialogue turns from the WIRED dataset. We use the previous metrics (W&A BERT-seq) found by Wachsmuth and Alshomary (2022) as our baseline. Percentages under each of the acts show micro- $F_1$  scores.

Dialogue acts	D01	D02	D03	D04	D05	D06	D07	D08	D09	Macro- $F_1$	Span Al.
BERT	73.14 %	72.72 %	74.02 %	55.43 %	50.25 %	66.28 %	<b>60.59</b> %	43.14 %	<b>94.86</b> %	69.01 %	-
SB2	21.66 %	54.27 %	2.83 %	7.63 %	39.16 %	9.03 %	33.66 %	22.78 %	93.50 %	28.72 %	59.61 %
GPT-3.5	19.71 %	54.73 %	11.69 %	0.00~%	8.70 %	7.01 %	19.74 %	12.98 %	83.87 %	22.30 %	59.41 %
GPT-4	34.25 %	57.44 %	2.38 %	14.32 %	27.13 %	9.42 %	34.71 %	28.97~%	93.12 %	31.45 %	63.18 %
Expl. acts	E01	E02	E03	E04	E05	E06	E07	E08	E09	Macro- $F_1$	Span Al.
BERT	64.66 %	<b>67.21</b> %	<b>94.69</b> %	<b>72.81</b> %	<b>64.80</b> %	<b>69.09</b> %	<b>64.99</b> %	80.65 %	80.34 %	75.89 %	-
SB2	8.93 %	33.63 %	89.08 %	56.00 %	31.67 %	17.97 %	20.21 %	0.00~%	4.64 %	26.22 %	60.54 %
GPT-3.5	20.06 %	10.02 %	84.27 %	24.23 %	16.90 %	19.35 %	4.69 %	0.00~%	7.07~%	18.66 %	49.72 %
GPT-4	27.70 %	42.11 %	86.18 %	66.52 %	34.82 %	42.93 %	19.94 %	9.07 %	20.77~%	35.00 %	61.49 %
Teaching acts	T01	Т02	Т03	Т04	Т05	T06	Т07	Т08	Т09	Macro- $F_1$	Span Al.
BERT	81.57 %	62.38 %	85.00 %	80.85 %	<b>89.61</b> %	86.34 %	85.67 %	<b>79.57</b> %	<b>72.91</b> %	82.36 %	-
SB2	28.24 %	28.04 %	13.23 %	8.42 %	49.12 %	7.83 %	2.09 %	10.21 %	29.44 %	19.62 %	44.39 %
GPT-3.5	22.89 %	8.95 %	19.10 %	7.25 %	40.31 %	10.31 %	11.80 %	5.13 %	13.66 %	15.49 %	31.55 %
GPT-4	31.21 %	26.32 %	16.80 %	5.58 %	47.23 %	14.71 %	16.02 %	10.93 %	25.63 %	21.60 %	39.55 %

Table 4: Language models evaluated on the tasks of sequence-labeling dialogue, explanation and teaching acts within dialogue turns from our ReWIRED dataset. Percentages under each of the acts show micro- $F_1$  scores. Act 10 was disregarded due to low number of instances, close-to-zero scores and irrelevance for the overall performance.

Still, they can handle the majority classes reasonably well (D02, E04, T05) or very well (D09, E03). However, in all other cases, all LLMs fail to assign the correct label consistently enough. Between the models, GPT-4 has a slight edge over SB2, which in turn is a lot more accurate than GPT-3.5. The difference in model performance is more pronounced for the already established acts (dialogue, explanation), but less so for our new teaching acts, whose label taxonomy is unlikely part of their training data. Evaluating how well the extracted spans align with human-annotated spans (rightmost column) reveals a similar pattern, i.e. GPT-4 beating the rest, except SB2 coming out on top for the teaching acts.

333

334

335

337

339

341

342

343

345

The prompt design that elicits structured prediction in the form of JSON objects from LLMs causes major problems for post-processing. After rigorously handling edge cases, we still find that 12.82 % of SB2, 9.73% of GPT-3.5 and 3.18 % of GPT-4 outputs result in invalid, unparseable JSONs. While the first two models repeatedly predicted duplicate spans with the same labels, in rare cases with an alternative label, GPT-4 tends to continue with a rationale explaining the annotation. These findings reflect challenges reported by concurrent related work applying LLMs to dialogue-related tasks (Zhao et al., 2023) and span-labeling tasks (Ziems et al., 2024; Wang et al., 2023) and the general difficulty of applying them to teaching settings (Wang and Demszky, 2023; Macina et al., 2023). 356

357

358

359

360

361

362

363

364

365

366

369

370

371

372

373

BERT, on the other hand, easily outperformed the prompt-based LLMs across every single act. The stark difference can be attributed to the importance of fine-tuning and the constraint to predict one of the ten acts. For span-labeling tasks like this, especially in the teaching act dimension where the performance is the overall highest, we recommend practitioners to look into such a controlled setup instead of unreliable prompt designs.

## 5.3 Quality Events in Instructional Explanations

Based on our annotation schema and as an addi-<br/>tional analysis, we develop and propose a test suite374based on didactics research. This novel assess-<br/>ment framework, which is termed as IXQUISITE,376

Category	Description	Origin	Measure
Check for prior knowledge	The teacher inquires the student about prior knowledge, back- ground, or what their interests might be	Kulgemeyer and Schecker (2009), Leinhardt and Steele (2005)	T01
Mindfulness of com- mon misconceptions	The teacher addresses common misconceptions	Wittwer et al. (2010), Andrews et al. (2011)	T04
Rule-Example struc- ture	The teacher states the abstract form of the concept being taught. Then the teacher gives some example to assist understanding	Tomlinson and Hunt (1971)	$T05 \rightarrow T03$
Example-Rule struc- ture	For procedural knowledge, the teacher first provides examples and then derives the general rule from them	Champagne et al. (1982)	$T03 \rightarrow T05$
Example/Analogy connection	The teacher explains how parts of the analogy/example relate to the concept being explored	Ogborn et al. (1996), Valle and Callanan (2006)	T06
Check for under- standing	The teacher tests the understanding of the student	Webb et al. (1995)	T08; E01
Remedial explana- tions	Either the teacher praises correct understanding (positive rein- forcement) or corrects improper understanding	Roelle et al. (2014), Sánchez et al. (2009)	E08

Table 5: Explanation and teaching acts-related measures in IXQUISITE for instructional explanation quality based on occurrences of classes from our annotation schema.

Category	Description	Origin	Measure
Minimal explana- tions	Low cognitive load, e.g. avoid redundancies (verbosity) such as introducing named entities	Black et al. (1986)	Frequency of named entities
Lexical complex- ity	The level of difficulty associated with any given word form by a particular individual or group	Kim et al. (2016)	Frequency of difficult words
Synonym density	Children are proven better aligned with consistent terminology; experts allow more synonyms	Wittwer and Ihme (2014)	Frequency of synonyms for the $n$ terms most connected to the topic
Correlation to teaching model	Correlation of teaching act order to prescribed teaching models	Oser and Baeriswyl (2002), Krabbe et al. (2015)	Edit distance between T01-T08 (asc.) and actual occurrences
Adaptation	The teacher incorporates prior knowledge, miscon- ceptions and interests and uses analogies	Wittwer et al. (2010)	Inverse frequency of synonyms in the text
Readability level	Indicator of how difficult a passage is to understand	Crossley et al. (2017)	Flesch-Kincaid Grade level
Coherence	How sentences relate to each other to create a logical and meaningful flow for the reader or listener	Lehman and Schraw (2002), Duffy et al. (1986)	Frequency of conjunctions and linking language

Table 6: Categories for instructional explanation quality and associated numerical measures in IXQUISITE.

addresses both the *form* of instructional explanations (in terms of syntax, vocabulary, etc) and their *function* (as present in the form of different classes in our annotation). While we only carry out analyses on and evaluate the **ReWIRED** dataset, we are confident that IXQUISITE can be applied to other kinds of instructional explanations, both humanand LLM-generated, among others.

378

380

381

384

The IXQUISITE test suite Since teaching models propose themselves as a proper method for instantiating learning, evaluating teaching according to their adherence to the prescribed method is also natural. We find that teaching models can serve as a quality metric and an opportunity to operationalize many other proposed evaluation metrics from didactics. We provide a new way to interact with the problem by providing a suite of tools that measure quality based on a large selection of proposed quality features from didactics literature. Through our suite of low-level quality tests, we aim to verify didactics theory in a controlled environment at a relatively low cost (using existing libraries, e.g., NLTK, SPACY, and TEXTSTAT). Following the literature review by Kulgemeyer (2018), we were able to track a list of seven events, which, when detected, have been shown to correlate to better learning outcomes, and seven more numerical metrics, which are the discrete values resulting from properties associated with better learning outcome. The events and metrics, along with their descriptions, are listed in Table 5 and Table 6, respectively. 396

397

398

399

400

401

402

403

404

405

406

407

408

IXQUISITE resultsThe qualitative act-based409measures, as well as the metrics correlate well with410the expert levels present in the ReWIRED dataset,411as seen in Figure 5. In terms of the former, test-412ing for understanding and remedial explanations413



Figure 5: IXQUISITE results with scores from explanation and teaching act-related measures (Table 5; top) and for the five levels in ReWIRED by category according to Table 6 (bottom).

are mostly present in lower expertise levels, which 414 is expected. Mindfulness (of common misconcep-415 tions) is especially high for colleague-level expla-416 nations and reflects the variation in conversation 417 topics present in the dataset. It is also interesting 418 to note that both *rule-example* and *example-rule* 419 structures are exceptionally present as well as in 420 teenager- and colleague-addressed dialogues. 421

422 Regarding our numerical metrics, we observe that explanations tailored to a child present a lower 423 bound across all our metrics, including a lower 424 lexical complexity, reading grade, synonym den-425 sity, and coherence. However, a general trend is 426 that graduate-level explanations score higher than 427 colleague-grade explanations (e.g., teaching model 428 correlation), probably because they are more fo-429 cused towards the actual topic of discussion, while 430 colleague-grade dialogues might also contain chit-431 chat and other topics, thus not necessarily follow-432 ing a teaching-like approach. In the case of adap-433 tation, graduate-level explanations are an outlier, 434 435 where the score is lower, which is a surprising result. Lastly, minimal explanations' scores for chil-436 dren average higher, possibly because of an attempt 437 to establish a common ground with world knowl-438 edge via entities. 439

### 6 Conclusion

We presented an extended dataset of instructional explanation dialogues in one-to-one tutorial sessions, **ReWIRED**, adding span-level annotations and new teaching acts dimension reflecting good practices according to didactics. Our language model analyses on the span-labeling tasks show that LLMs, including GPT-4, fall far behind controlled setups like a fine-tuned BERT in reliably detecting acts across multiple act dimensions. Our IXQUISITE suite of metrics for quality events in instructional explanations represent the different expertise levels of explainees well and are a first step in operationalizing pedagogical psychological theory for tutorial dialogues in NLP. 440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

In the future, we plan to follow concurrent work in exploring LLM-based explanation quality evaluation (Rooein et al., 2023), especially for metrics such as Adaptation and Mindfulness of common misconceptions. These are hard to capture with the more traditional approach we chose and instead require world knowledge that LLMs can provide. Further data collection and fine-tuning will also allow mimicking the behavior found in classroom transcripts for multi-turn systems. This forms a fertile basis for more satisfactory explanation dialogues from automated tutoring systems.

568

569

570

571

### 467 Limitations

477

478

479

480

481

482

483

484

485

486

487

488

489

490

492

493

494

495

496

497

498

499

500

502

504

505

507

509

510

511

512

513

514

Resulting from the low inter-annotator agreement 468 for the teaching acts as discussed in §3.2, we want 469 to perform data collection involving teachers and 470 didacticians in the future. Additionally, a portion of 471 our test suite relies on human annotation, a factor 472 that may introduce inconsistencies. In this case, 473 replication or extension of the test suite might be 474 difficult without a reliable teaching act prediction 475 model. 476

> Due to time and budget constraints, we were not able to explore many different prompt patterns in our LLM experiments. The prompt design utilized in our study may not represent an ideal formulation, potentially influencing the model's performance.

The dataset we present is extracted from videos in the transcription, audio and visual elements are not present. The efficacy of our approach may vary depending on the complexity and diversity of the multimodal inputs, if present.

Last but not least, the generalizability of our findings may be constrained by the narrow domain of dialogues examined, limiting extrapolation to broader conversational contexts.

### 491 Ethical statement

We do not see any immediate ethical concerns with respect to research and development. The data included in the corpus is readily available from the WIRED web resources. The nine annotators in our study were paid the minimum wage in conformance with the standards of our host institutions' regions. In our view, the provided prediction models target dimensions of dialogue turns that are not prone to be misused for ethically doubtful applications.

### References

- Elaine Allensworth, Macarena Correa, and Steve Ponisciak. 2008. From high school to the future: Act preparation-too much, too late. why act scores are low in chicago and what it means for schools. *Consortium on Chicago School Research*.
- Tessa M Andrews, Michael J Leonard, Clinton A Colgrove, and Steven T Kalinowski. 2011. Active learning not associated with student learning in a random sample of college biology courses. *CBE—Life Sciences Education*, 10(4):394–405.
- Deborah Loewenberg Ball and Brian Rowan. 2004. Introduction: Measuring instruction. *The Elementary School Journal*, 105(1):3–10.

- Emily M. Bender and Alexander Koller. 2020. Climbing towards NLU: On meaning, form, and understanding in the age of data. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5185–5198, Online. Association for Computational Linguistics.
- John B. Black, John M. Carroll, and Stuart M. McGuigan. 1986. What kind of minimal instruction manual is the most effective. *SIGCHI Bull.*, 18(4):159–162.
- Melissa Boston. 2012. Assessing instructional quality in mathematics. *The Elementary School Journal*, 113(1):76–104.
- Sarah Bourse and Patrick Saint-Dizier. 2012. A repository of rules and lexical resources for discourse structure analysis: the case of explanation structures. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 2778–2785, Istanbul, Turkey. European Language Resources Association (ELRA).
- Harry Bunt, Jan Alexandersson, Jae-Woong Choe, Alex Chengyu Fang, Koiti Hasida, Volha Petukhova, Andrei Popescu-Belis, and David Traum. 2012. ISO 24617-2: A semantically-based standard for dialogue annotation. In Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12), pages 430–437, Istanbul, Turkey. European Language Resources Association (ELRA).
- Andrew Caines, Helen Yannakoudakis, Helen Allen, Pascual Pérez-Paredes, Bill Byrne, and Paula Buttery. 2022. The teacher-student chatroom corpus version 2: more lessons, new annotation, automatic detection of sequence shifts. In *Proceedings of the 11th Workshop on NLP for Computer Assisted Language Learning*, pages 23–35, Louvain-la-Neuve, Belgium. LiU Electronic Press.
- Audrey B Champagne, Leopold E Klopfer, and Richard F Gunstone. 1982. Cognitive research and the design of science instruction. *Educational Psychologist*, 17(1):31–53.
- Mark G. Core and James F. Allen. 1997. Coding dialogs with the DAMSL annotation scheme. In AAAI fall symposium on communicative action in humans and machines, volume 56, pages 28–35. Boston, MA.
- Scott A. Crossley, Stephen Skalicky, Mihai Dascalu, Danielle S. McNamara, and Kristopher Kyle. 2017. Predicting text comprehension, processing, and familiarity in adult readers: New approaches to readability formulas. *Discourse Processes*, 54(5-6):340–359.
- Jorge Del-Bosque-Trevino, Julian Hough, and Matthew Purver. 2021. Communicative grounding of analogical explanations in dialogue: A corpus study of conversational management acts and statistical sequence models for tutoring through analogy. In *Proceedings of the Reasoning and Interaction Conference* (*ReInAct 2021*), pages 23–31, Gothenburg, Sweden. Association for Computational Linguistics.

Dorottya Demszky and Heather Hill. 2023. The NCTE transcripts: A dataset of elementary math classroom transcripts. In *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pages 528–538, Toronto, Canada. Association for Computational Linguistics.

572

574

582

584

585

587

591

594

595

596

597

598

612

614

615

616

617

618

619

622

627

- Dorottya Demszky, Jing Liu, Zid Mancenido, Julie Cohen, Heather Hill, Dan Jurafsky, and Tatsunori Hashimoto. 2021. Measuring conversational uptake: A case study on student-teacher interactions. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 1638–1653, Online. Association for Computational Linguistics.
- Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. 2019. Wizard of Wikipedia: Knowledge-powered conversational agents. In *International Conference on Learning Representations*.
- Gerald G. Duffy, Laura R. Roehler, Michael S. Meloth, and Linda G. Vavrus. 1986. Conceptualizing instructional explanation. *Teaching and Teacher Education*, 2(3):197–214.
- Peter Jansen, Elizabeth Wainwright, Steven Marmorstein, and Clayton Morrison. 2018. WorldTree:
   A corpus of explanation graphs for elementary science questions supporting multi-hop inference. In Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018), Miyazaki, Japan. European Language Resources Association (ELRA).
- Yea-Seul Kim, Jessica Hullman, Matthew Burgess, and Eytan Adar. 2016. SimpleScience: Lexical simplification of scientific terminology. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1066–1071, Austin, Texas. Association for Computational Linguistics.
- Heiko Krabbe, Simon Zander, and Hans Ernst Fischer. 2015. Lernprozessorientierte Gestaltung von Physikunterricht - Materialien zur Lehrerfortbildung. Waxmann.
- Christoph Kulgemeyer. 2018. Towards a framework for effective instructional explanations in science teaching. *Studies in Science Education*, 54(2):109–139.
- Christoph Kulgemeyer and Horst Schecker. 2009. Kommunikationskompetenz in der physik: Zur entwicklung eines domänenspezifischen kompetenzbegriffs. Zeitschrift für Didaktik der Naturwissenschaften, 15:131–153.
- Ashlee Kupor, Candice Morgan, and Dorottya Demszky. 2023. Measuring five accountable talk moves to improve instruction at scale. *arXiv*, abs/2311.10749.
- Yoonjoo Lee, Tae Soo Kim, Sungdong Kim, Yohan Yun, and Juho Kim. 2023. DAPIE: Interactive stepby-step explanatory dialogues to answer children's

why and how questions. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, CHI '23, New York, NY, USA. Association for Computing Machinery. 628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

- Stephen Lehman and Gregory Schraw. 2002. Effects of coherence and relevance on shallow and deep text processing. *Journal of Educational Psychology*, 94(4):738–750.
- Gaea Leinhardt and Michael D. Steele. 2005. Seeing the complexity of standing to the side: Instructional dialogues. *Cognition and Instruction*, 23(1):87–163.
- Jakub Macina, Nico Daheim, Lingzhi Wang, Tanmay Sinha, Manu Kapur, Iryna Gurevych, and Mrinmaya Sachan. 2023. Opportunities and challenges in neural dialog tutoring. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2357–2372, Dubrovnik, Croatia. Association for Computational Linguistics.
- Dakota Mahan, Ryan Carlow, Louis Castricato, Nathan Cooper, and Christian Laforte. 2023. Stable Beluga models. Hugging Face.
- Lindsay Clare Matsumura, Helen E. Garnier, Sharon Cadman Slater, and Melissa D. Boston. 2008. Toward measuring instructional interactions "at-scale". *Educational Assessment*, 13(4):267–300.
- Danielle S. McNamara, Arthur C. Graesser, Philip M. McCarthy, and Zhiqiang Cai. 2014. *Automated Evaluation of Text and Discourse with Coh-Metrix*. Cambridge University Press.
- Tim Miller. 2019. Explanation in artificial intelligence: Insights from the social sciences. *Artificial Intelligence*, 267:1–38.
- Nasrin Mostafazadeh, Aditya Kalyanpur, Lori Moon, David Buchanan, Lauren Berkowitz, Or Biran, and Jennifer Chu-Carroll. 2020. GLUCOSE: GeneraLized and COntextualized story explanations. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 4569–4586, Online. Association for Computational Linguistics.
- Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. 2023. Orca: Progressive learning from complex explanation traces of GPT-4. *arXiv*, abs/2306.02707.
- Hiroki Nakayama, Takahiro Kubo, Junya Kamura, Yasufumi Taniguchi, and Xu Liang. 2018. doccano: Text annotation tool for human. Software available from https://github.com/doccano/doccano.
- Jon Ogborn, Gunther Kress, Isabel Martins, and Kieran McGillicuddy. 1996. *Explaining science in the classroom*. McGraw-Hill Education (UK).

Fritz Oser and Franz Baeriswyl. 2002. AERA's Handbook of Research on Teaching, 4th Edition, pages 1031–1065. Washington: American Educational Research Association (AERA).

684

695

697

704

710

711

712

713

714

715

716

717

718

719

720

726

727 728

729

730

731

732

733

734

736

- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics.
- Julian Roelle, Kirsten Berthold, and Alexander Renkl. 2014. Two instructional aids to optimise processing and learning from instructional explanations. *Instructional Science*, 42:207–228.
- Julian Roelle, Claudia Müller, Detlev Roelle, and Kirsten Berthold. 2015. Learning from instructional explanations: Effects of prompts based on the activeconstructive-interactive framework. *PLOS ONE*, 10(4):e0124115.
- Donya Rooein, Amanda Cercas Curry, and Dirk Hovy. 2023. Know your audience: Do LLMs adapt to different age and education levels? *arXiv*, abs/2312.02065.
- Emilio Sánchez, Héctor García-Rodicio, and Santiago R Acuna. 2009. Are instructional explanations more effective in the context of an impasse? *Instructional Science*, 37:537–563.
- Hendrik Schuff, Heike Adel, Peng Qi, and Ngoc Thang Vu. 2023. Challenges in explanation quality evaluation. *arXiv*, abs/2210.07126v2.
- Katherine Stasaski, Kimberly Kao, and Marti A. Hearst. 2020. CIMA: A large open access dialogue dataset for tutoring. In *Proceedings of the Fifteenth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 52–64, Seattle, WA, USA → Online. Association for Computational Linguistics.
- Andreas Stolcke, Klaus Ries, Noah Coccaro, Elizabeth Shriberg, Rebecca Bates, Daniel Jurafsky, Paul Taylor, Rachel Martin, Carol Van Ess-Dykema, and Marie Meteer. 2000. Dialogue act modeling for automatic tagging and recognition of conversational speech. *Computational Linguistics*, 26(3):339–374.
- Peter D Tomlinson and David E Hunt. 1971. Differential effects of rule-example order as a function of learner conceptual level. *Canadian Journal of Behavioural Science/Revue canadienne des sciences du comportement*, 3(3):237.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura,

Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and finetuned chat models. *arXiv*, abs/2307.09288.

737

738

740

741

742

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

778

779

780

781

782

783

784

785

787

788

789

790

791

792

793

794

- Araceli Valle and Maureen A Callanan. 2006. Similarity comparisons and relational analogies in parent-child conversations about science topics. *Merrill-Palmer Quarterly* (1982-), pages 96–124.
- Henning Wachsmuth and Milad Alshomary. 2022. "mama always had a way of explaining things so I could understand": A dialogue corpus for learning to construct explanations. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 344–354, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Rose Wang and Dorottya Demszky. 2023. Is ChatGPT a good teacher coach? measuring zero-shot performance for scoring and providing actionable insights on classroom instruction. In *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pages 626–667, Toronto, Canada. Association for Computational Linguistics.
- Shuhe Wang, Xiaofei Sun, Xiaoya Li, Rongbin Ouyang, Fei Wu, Tianwei Zhang, Jiwei Li, and Guoyin Wang. 2023. GPT-NER: Named entity recognition via large language models. *arXiv*, abs/2304.10428.
- Noreen M Webb, Jonathan D Troper, and Randy Fall. 1995. Constructive activity and learning in collaborative small groups. *Journal of educational psychology*, 87(3):406.
- Jacob Whitehill and Jennifer LoCasale-Crouch. 2024. Automated evaluation of classroom instructional support with LLMs and BoWs: Connecting global predictions to specific feedback. *arXiv*, abs/2310.01132v2.
- Jörg Wittwer, Matthias Nückles, Nina Landmann, and Alexander Renkl. 2010. Can tutors be supported in giving effective explanations? *Journal of Educational Psychology*, 102(1):74.
- Jörg Wittwer and Natalie Ihme. 2014. Reading skill moderates the impact of semantic similarity and causal specificity on the coherence of explanations. *Discourse Processes*, 51(1-2):143–166.
- Weixiang Zhao, Yanyan Zhao, Xin Lu, Shilong Wang, Yanpeng Tong, and Bing Qin. 2023. Is Chat-GPT equipped with emotional dialogue capabilities? *arXiv*, abs/2304.09582.

Caleb Ziems, William Held, Omar Shaikh, Jiaao Chen,
Zhehao Zhang, and Diyi Yang. 2024. Can large language models transform computational social science? *Computational linguistics*, 50(1).

### 799 Appendix

800

804

805

806

807

813

814

815

816

817

### A Examples for acts

Figures 6, 7, and 8 show examples from ReWIRED for each of the acts as provided to the annotators.

### **B** Label distributions

Figure 9 shows the distribution of annotated acts in the dialogue and explanation dimensions.

### C Models

808	• MiniLM	sentence	embeddings:	https://
809	hugging	face.co/s	entence-trans	sformers/
810	all-Min	iLM-L6-v2	2	
811	• BERT:		https://huggi	ngface.co/
812	bert-ba	ise-uncase	ed	

 Stable Beluga 2: https://huggingface.co/ petals-team/StableBeluga2

• GPT-3.5 and GPT-4: https://platform. openai.com/docs/api-reference/chat

### D Prompt design

Figure 10 and Figure 11 depict the prompts used
with SB2, GPT-3.5 and GPT-4 to produce the predictions whose evaluation is shown in Table 3 and
Table 4, respectively.

Explainer: Do you like science? D01 - Chec Select a label D01 - Check Question 😵	Child: What's blockchain? •D02 - Wh Select a label D02 - What/How Question (2)
(a) D01: Check Question	(b) D02: What/How Question
Child: It's when you give up most of what you want, right? •D01 - Check Explainer: When you give up most of what you want? •D03 - Other Explainer: Well, sometimes that definitely happens for sure.	Child: <u>Really?</u> •D01 - Check Explainer: <u>Yeah.</u> •DC <sup>4</sup> Select a label
•D07 - Agreei	D04 - Confirming Answer 😵
(c) D03: Other Question	(d) D04: Confirming Answer
Teenager: [laughing] Yeah. *D06 - Other Select a label D06 - Other Answer ()	Child: well, I would trade something I don't like so much. Explainer: That's probably a good idea, maybe somebody else • D07 - Agreei Select a label 
(e) Doo: Other Allswer	(1) D07: Agreeing Statement
	•D02 - What/H Teenager: I like biology and computer science. •D09 - Inform
Explainer: It's actually a way that we can trade.	Explainer: All right, you're in the wrong place. •D10 - Anothe Select a label D10-Another Act •D09 - Info (h) D10: Other Act

Figure 6: Examples for dialogue acts. D05 and D08 are left out, because they are analogous to D04 and D07, respectively.

Explainer: So based on what we've discussed today, • E01 - Test U		Explainer So have	vou takan a quantum machanica courae?		
Explainer: in your own words, what is a zero-knowledge	proof	• E02 - T	e Select a label E02. Test Princ Knowledge	•	<b>^</b>
(a) E01: Test Understanding	proof.	(	(b) E02: Test Prior Knowledge		
Explainer: What it teaches us					
•E04 - Reques					
Explainer: as we make these devices smaller and	smaller,				
Explainer: their properties begin to now depend		Explainer: What ma •E04 - Re	ade you choose that? eques Select a label		•
Explainer: on the size and the orientation of these	e devices.	Undergrad: Like an • E03 - F	y fi 1 E01 - Test Understanding		
(c) E03: Provide Explanation (The color/label is w	vrong here!)		(d) E04: Request Explanation		
Explainer: So based on what we've discussed today, •E01 - Test U					
Explainer: in your own words, what is a zero-knowled	dge proof.				
Teenager: It's like, if you have this really important so	ecret				
•E05 - Signal		Explainer: it's a qu •E03 - Provid	uantum computer		
Teenager: that you want somebody to know about,		Toopogor: A what	2		
Teenager: but you don't want to tell them everything.	i	•E06 - Signal			
(e) E05: Signal Understanding		(f)	E06: Signal Non-understanding		
	Explainer: But wh •E08 -	nat if I could prove to Provid	у уоц		
	Explainer: that I k	know where the puff	in is		
	Explainer: withou	it revealing to you w	here it is?		
	Explainer: Let me	e show you.			
	Explainer: I took	that photo that we s	howed you.		
Teenager: I do, I try. [laughs]	Explainer: And I p	out it behind this pos	ster here.		
•E10 - Other	Explainer: Why do	on't you go take a lo	ok through that hole?		
Explainer: <u>You try, we all gotta try.</u> •E07 - Provid	Child: I see the pr •E05 - Sign	uffin. nal			
(g) E07: Provide Feedback		(h) E03	8: Provide Assessment		
Explainer: So what's your major?					
		Te	enager: Wow.		
Undergrad: Chemical engineering. • E09 - Provid Select a label		۰E	10 - Other		
E09 - Provide extraneous information	- Ex	plainer: so my team is working on buildi	ng		
(i) E09: Provide Extraneous Info	rmation		(j) E10: Other Act		

Figure 7: Examples for Explanation Acts.



Figure 8: Examples for teaching acts T01-T06 and T09. Examples for T07 and T08 are in Figure 3.





(a) Distribution of dialogue acts Figure 9: Distribution of annotated acts in ReWIRED across the five expertise levels for three dimensions dialogue (a) and explanation (b).

```
1 system_prompt = (f"You are an expert annotator. In the following, you will be requested to
   \rightarrow classify a single turn of a dialogue between explainer and {student_role}.\n")
  # Example label mapping (dialogue acts)
2
   WIRED_da_label_mapping = {
3
        (D01) To ask a check question': 1,
4
       '(D02) To ask what/how question': 2,
5
       '(D03) To ask other kind of questions': 3,
6
       '(D04) To answer a question by confirming': 4,
7
8
       '(D05) To answer a question by disconfirming': 5,
       '(D06) To answer - Other': 6,
9
10
       '(D07) To provide agreement statement': 7,
11
       '(D08) To provide disagreement statement': 8,
       '(D09) To provide informing statement': 9,
12
       '(D10) Other': 10,
13
14 }
15 label_schema = ("The label schema consists of the following 10 classes:\n* " + "\n*

. join(list(WIRED_da_label_mapping.keys())) + "\n")

16
   read_instruction = f"The excerpt from the dialogue:\n{turn_text}\n"
17 task_instruction = "Predicted label:\n"
18 # Combine inputs to single string
entire_prompt = system_prompt + label_schema + read_instruction + task_instruction
```

Figure 10: Simplified version of the Python code showing the turn classification task prompt for WIRED.

```
system_prompt = (f"You are an expert annotator. ")
read_instruction = (f"Here is one turn from a dialogue between an explainer and a {student_role}
on the topic of {topic}:\n{turn_text}\n")
task_instruction = ("Please extract the spans from the turn and assign a label to each of the
spans. It is possible that the whole turn is just one span, because the act applies to its
entirety. Please present your predictions in a JSON format like this: {\n\t{\n\t\'Span':
} '...', \n\t\t'Predicted label': '...' \n\t},\n]\n")
entire_input = system_prompt + read_instruction + label_schema + task_instruction
```

Figure 11: Simplified version of the Python code showing the span labeling task prompt for ReWIRED.